# Just call it a "*p*-value"!
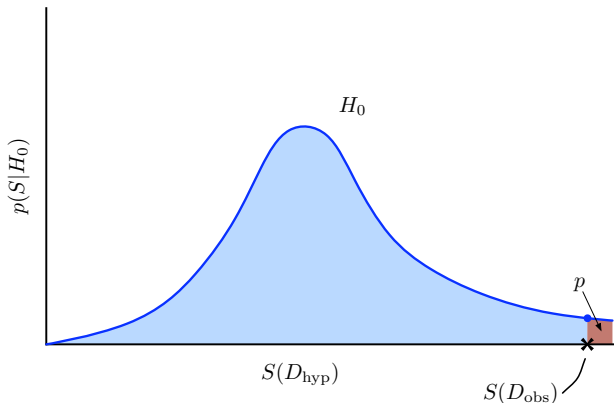
*(not a hypothesis probability,*
*not a false alarm probability)*

Tom Loredo
Cornell Center for Astrophysics and Planetary Science
http://www.astro.cornell.edu/staff/loredo/

CASt Summer School — 2–3 June 2016

# *p*-**values**

$$p = P(S(D) > S(D_{\mathrm{obs}})|H_0)$$



Smaller *p*-values indicate stronger evidence against $H_0$.

Astronomers call this the *significance level* or (sometimes) *false-alarm probability*. Statisticians don't—for good reason!

# An old misunderstanding

## THE ABSOLUTE MAGNITUDE DISTRIBUTION OF KUIPER BELT OBJECTS

### ABSTRACT

Here we measure the absolute magnitude distributions ($H$-distribution) of the dynamically excited and quiescent (hot and cold) Kuiper Belt objects (KBOs), and test if they share the same $H$-distribution as the Jupiter Trojans. From

"The Kolmogorov–Smirnov test reveals that *the probability that the Trojans and cold KBOs share the same parent H-distribution is less than 1 in 1000*."

A coauthor collaborated with me on earlier ~~astrostat~~ papers—*ouch!*

$p$-values are probabilities for *data*, not hypotheses—only Bayesian methods can give probabilities for hypotheses

# A newer misunderstanding

"This detection has a signal-to-noise ratio of 4.1 with an empirically estimated upper limit on false alarm probability of 1.0%."

"...the false alarm probability for this signal is rather high at a few percent."

"This signal has a false alarm probability of <4 % and is consistent with a planet of minimum mass 2.2 $M_\odot$..."

"We find a false-alarm probability <10-4 that the RV oscillations attributed to CoRoT-7b and CoRoT-7c are spurious effects of noise and activity."

# **What's wrong?**

"*This signal, with $S(D)=S_{obs}$, has a FAP of p…*"

$p$ is not a property of **this** signal; rather, it's the size of the **ensemble** of possible null-generated signals with $S(D)>S_{obs}$

*Every one* of those signals is a false alarm: each one has a FAP=1 in the context producing the $p$-value!

For any signal to have FAP≠1, alternatives to the null must sometimes act; the FAP will depend on how often they do (and what they are)

What a *p*-value really means:

(In the voice of Don LaFontaine or Lake Bell)

⊞

"*In a world… with absolutely no sources, with a threshold set so we wrongly claim to detect sources 100×p% of the time, this data would be judged a detection—and it would be the data providing the weakest evidence for a source in that world.*"
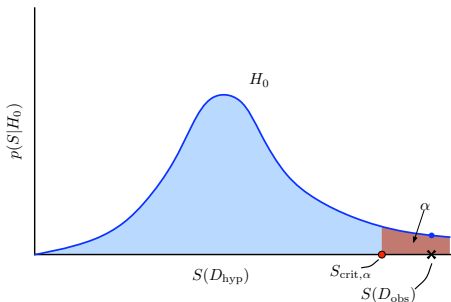
Who wants to say that?!

Whence "*p-value*" — a measure of "surprisingness" under the null whose main virtue is that p is uniformly distributed under the null

# Significance Testing and $p$-values

*Neyman-Pearson testing*

- Specify simple null hypothesis $H_0$ such that rejecting it implies an interesting effect is present
- Devise statistic $S(D)$ measuring departure from null
- Divide sample space into probable and improbable parts (for $H_0$); $p(\text{improbable}|H_0) = \alpha$ (Type I error rate), with $\alpha$ specified a priori
- If $S(D_{\text{obs}})$ lies in improbable region, reject $H_0$; otherwise accept it
- Report: "$H_0$ was rejected (or not) with a procedure with false-alarm frequency $\alpha$"

Neyman and Pearson devised this approach guided by
Neyman's *frequentist principle*:

> *In repeated practical use of a statistical procedure, the
> long-run average actual error should not be greater than
> (and ideally should equal) the long-run average reported
> error. (Berger 2003)*

A *confidence region* is an example of a familiar procedure
satisfying the frequentist principle.

They insisted that one also specify an alternative, and find the
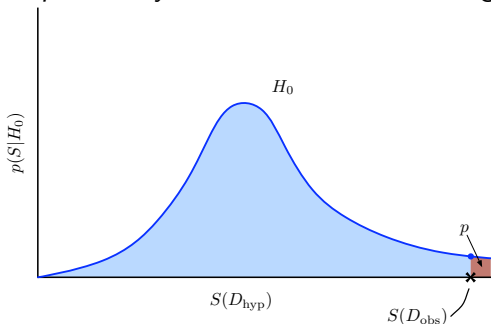error rate for falsely rejecting it (Type II error).

## Fisher's p-value testing

Fisher (and others) felt reporting a rejection frequency of $\alpha$ no matter where $S(D_{\mathrm{obs}})$ lies in the rejection region does not accurately communicate the strength of evidence against $H_0$.

He advocated reporting the *p-value*:

$$p = P(S(D) > S(D_{\mathrm{obs}})|H_0)$$

Smaller *p*-values indicate stronger evidence against $H_0$.

Astronomers call this the *significance level* or (sometimes) *false-alarm probability*. Statisticians don't—for good reason!

### p-value complications

Fisherian testing does not have the straightforward frequentist properties of NP testing, but everyone uses it anyway.

E.g., rejections of $H_0$ with $p$-value$= 0.05$ are *not* "wrong 5% of the time under the null" or "with 5% false-alarm probability." They are wrong *100% of the time* under the null. To quantify the conditional error rate (i.e., the error rate among datasets with the same $p$-value), you *must* say something about the alternative.

Even NP tests have unpleasant frequentist properties; e.g., the strength of the evidence against the null (e.g., quantified by a conditional false alarm rate) for a fixed-$\alpha$ test grows weaker as $N$ increases. NP themselves advocated decreasing $\alpha$ with $N$, but there are no general rules for this.

## *False alarm rates*

Berger (2003) discusses the relationship between *p*-values, false alarm rates, and Bayesian posterior probabilities (or odds and Bayes factors).

In simple settings where one can easily bound false alarm rates, he shows the *p*-value significantly underestimates the false alarm rate among datasets sharing a given *p*-value.

This gives insight into why we've come to consider apparently small *p*-values—like "$2\sigma$" ($p \approx 0.05$) or "$3\sigma$" ($p \approx 0.003$)—to represent only weak evidence against the null. Typically, datasets with such *p*-values are not much more probable under alternatives than under the null.

He also shows that a "conditional frequentist" calculation of the false alarm rate in some settings amounts to computation of a Bayes factor. (See example below.)

## Entries to the literature

- "402 Citations Questioning the Indiscriminate Use of Null Hypothesis Significance Tests in Observational Studies" (Thompson 2001) [web site]

- *The significance test controversy: a reader* (ed. Morrison & Henkel 1970, 2006) [Google Books]

- "Could Fisher, Jeffreys and Neyman Have Agreed on Testing?" (Berger 2003 with discussion; 2001 Fisher award Lecture), *Statistical Science*, **18**, 1–32 [journal site—highly recommended!]

- "Odds Are, It's Wrong: Science fails to face the shortcomings of statistics" (By Tom Siegfried 2010) [*Science News*, March 2010]

- "Scientific method: Statistical errors" (By Regina Nuzzo 2014) [*Nature* news feature, Feb 2014]

- "The ASA's statement on p-values: context, process, and purpose" [*The American Statistician*, March 2016]

## Example based on Berger (2003)

Model: $x_i = \mu + \epsilon_i$, $(i = 1 \text{ to } n)$ $\qquad \epsilon_i \sim N(0, \sigma^2)$

Null hypothesis, $H_0$: $\quad \mu = \mu_0 = 0$

Test statistic:

$$t(x) = \frac{|\bar{x}|}{\sigma/\sqrt{n}}$$

$p$-value:

$$
\begin{aligned}
p(t|H_0) &= \frac{1}{\sqrt{2\pi}} e^{-t^2/2} \\
p\text{-value} &= P(t \geq t_{\text{obs}})
\end{aligned}
$$

| $t$ | $p$-value |
| --- | --- |
| 1 | 0.317 |
| 2 | 0.046 |
| 3 | 0.003 |

$p = .05 \rightarrow$ "significant"

$p = .01 \rightarrow$ "highly significant"

Collect the *p*-values from a large number of tests in situations where the truth eventually became known, and determine how often $H_0$ is true at various *p*-value levels.

- Suppose that, overall, $H_0$ was true about half of the time.
- Focus on the subset with $t \approx 2$ (say, $[1.95, 2.05]$ so $p \in [.04, .05]$, so that $H_0$ was rejected at the 0.05 level.
- Find out how many times in that subset $H_0$ turned out to be true.
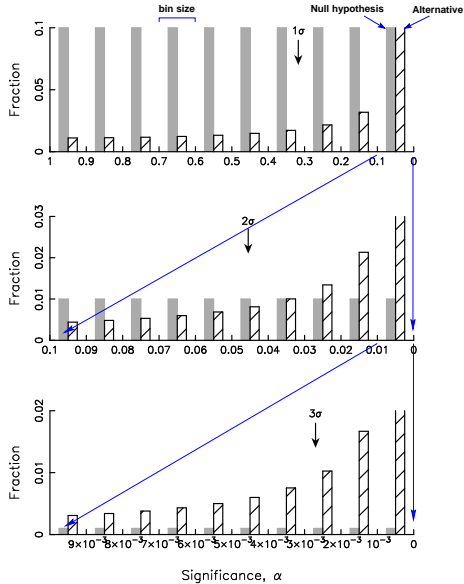- Do the same for other significance levels.

## A Monte Carlo experiment

- Choose $\mu = 0$ OR $\mu \sim N(0, 4\sigma^2)$ with a fair coin flip[*]

- Simulate $n$ data, $x_i \sim N(\mu, \sigma^2)$ (use $n = 20, 200, 2000$)

- Calculate $t_{\mathrm{obs}} = \frac{|\bar{x}|}{\sigma/\sqrt{n}}$ and $p(t_{\mathrm{obs}}) = P(t > t_{\mathrm{obs}}|\mu = 0)$

- Bin $p(t)$ separately for each hypothesis; repeat

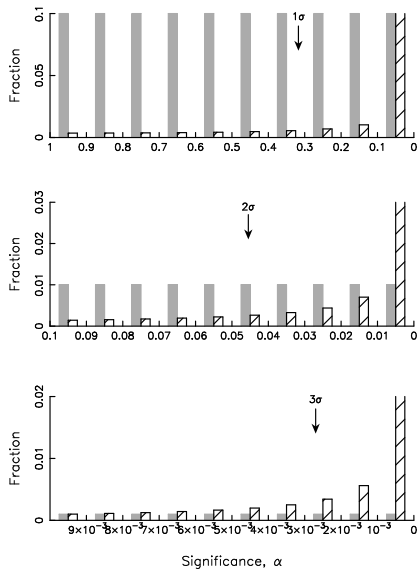Compare how often the two hypotheses produce data with a 1–, 2–, or 3–$\sigma$ effect.

[*]A neutral assumption that gives alternatives a "fair" chance and may *over*estimate the evidence against $H_0$ in real settings where the null is more prevalent
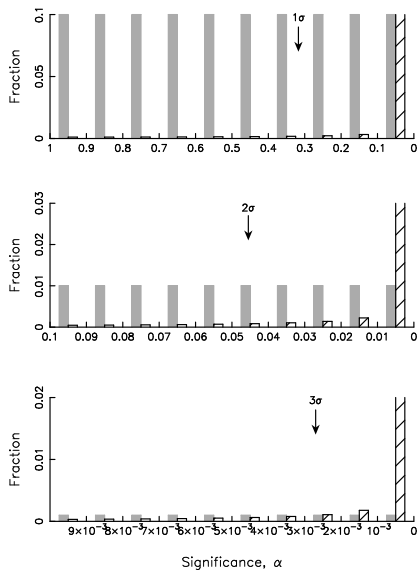
# Significance Level Frequencies, $n = 20$



Significance, $\alpha$

# Significance Level Frequencies, $n = 200$

# Significance Level Frequencies, $n = 2000$

What about another $\mu$ prior?

- For data sets with $H_0$ rejected at $p \approx 0.05$, $H_0$ will be true *at least* 23% of the time (and typically close to 50%). (Edwards et al. 1963; Berger and Selke 1987)
- At $p \approx 0.01$, $H_0$ will be true *at least* 7% of the time (and typically close to 15%).

What about a different "true" null frequency?

- If the null is initially true 90% of the time (as has been estimated in some disciplines), for data producing $p \approx 0.05$, the null is true at least 72% of the time, and typically over 90%.

In addition . . .

- At a fixed $p$, the proportion of the time $H_0$ is falsely rejected grows as $\sqrt{n}$. (Jeffreys 1939; Lindley 1957)
- Similar results hold generically; e.g., for $\chi^2$. (Delampady & Berger 1990)

*A p-value is not an easily interpretable measure of the weight of evidence against the null.*

- It does not measure how often the null will be wrongly rejected among similar data sets
- A naive false alarm interpretation typically overestimates the evidence
- For fixed *p*-value, the weight of the evidence decreases with increasing sample size

## Bayesian view of false-alarm rate

$$B \equiv \frac{p(\{x_i\}|H_1)}{p(\{x_i\}|H_0)} = \frac{p(p_{\mathrm{obs}}|H_1)}{p(p_{\mathrm{obs}}|H_0)}$$

$\rightarrow B$ here is just the ratio calculated in the Monte Carlo!

*Why is the p-value a poor measure of the weight of evidence?*

- We should be *comparing hypotheses*, not trying to identify rare/surprising events—an observation surprising under the null motivates rejection only if it is not surprising under reasonable alternatives

- Comparison should use the *actual data*, not merely membership of the data in some larger set. A *p*-value conditions on incomplete information.

Harold Jeffreys, addressing an audience of statisticians:

> For $n$ from about 10 to 500 the usual result is that $K = 1$ when $(a - \alpha_0)/s_\alpha$ is about 2... not far from the rough rule long known to astronomers, i.e. that differences up to twice the standard error usually disappear when more or better observations become available... I have always considered the arguments for the use of $P$ absurd. They amount to saying that a hypothesis that may or may not be true is rejected because a greater departure from the [observed] trial was improbable; that is, that it has not predicted something that has not happened. As an argument astronomer's experience is far better. (Jeffreys 1980)